

CHE384, From Data to Decisions: Measurement, Uncertainty, Analysis, and Modeling

## Lecture 48

### Standardized Variables

Chris A. Mack  
Adjunct Associate Professor

<http://www.lithoguru.com/scientist/statistics/>

© Chris Mack, 2016 Data to Decisions 1

## Standardizing Variables

- Consider the model  

$$y = f(x, \beta) + \varepsilon \quad x = (x_1, x_2, \dots); \quad \beta = (\beta_0, \beta_1, \dots)$$
- We can “standardize” each of the variables

	response	j <sup>th</sup> predictor
Standardized (unit normal scaling)	$\hat{y}_i = \frac{y_i - \bar{y}}{s_y}$	$\hat{x}_{ij} = \frac{x_{ij} - \bar{x}_j}{s_j}$
Correlation Transformation (unit length scaling)	$\hat{y}_i = \frac{y_i - \bar{y}}{\sqrt{n-1}s_y}$	$\hat{x}_{ij} = \frac{x_{ij} - \bar{x}_j}{\sqrt{n-1}s_j}$

© Chris Mack, 2016 Data to Decisions 2

## Recall OLS Matrix Formulation

$$Y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}, \quad \varepsilon = \begin{bmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{bmatrix}, \quad \beta = \begin{bmatrix} \beta_0 \\ \vdots \\ \beta_{p-1} \end{bmatrix}$$

$$X = \begin{bmatrix} 1 & x_{11} & \cdots & x_{1,p-1} \\ 1 & x_{21} & \cdots & x_{2,p-1} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \cdots & x_{n,p-1} \end{bmatrix} = \text{design matrix}$$

$$Y = X\beta + \varepsilon$$

(each row in X and Y is a “data point”)

© Chris Mack, 2016 Data to Decisions 3

## Recall OLS Matrix Formulation

$$Y = X\beta + \varepsilon$$

$$\hat{\beta} = (X^T X)^{-1} X^T Y$$

$$\hat{Y} = X\hat{\beta} = X(X^T X)^{-1} X^T Y = HY$$

$$H = X(X^T X)^{-1} X^T = \text{hat matrix}$$

$$\varepsilon = Y - \hat{Y} = (I - H)Y$$

$$\text{cov}(\hat{\beta}) = s_e^2 (X^T X)^{-1}$$

© Chris Mack, 2016 Data to Decisions 4

## Standardized Matrix Formulation

- The “standardized regression model” uses the correlation transformation
- Note: it will not have an intercept term  $\beta_0$

$$\tilde{Y} = \tilde{X}\tilde{\beta} + \varepsilon \quad \tilde{X} = \begin{bmatrix} \tilde{x}_{11} & \cdots & \tilde{x}_{1,p-1} \\ \vdots & \ddots & \vdots \\ \tilde{x}_{n1} & \cdots & \tilde{x}_{n,p-1} \end{bmatrix} = \text{design matrix}$$

Converting back:  $\beta_j = \tilde{\beta}_j \left( \frac{s_y}{s_j} \right) \quad \beta_0 = \bar{y} - \sum_{j=1}^{p-1} \beta_j \bar{x}_j$

$$SE(\beta_j) = SE(\tilde{\beta}_j) \left( \frac{s_y}{s_j} \right)$$

© Chris Mack, 2016 Data to Decisions 5

## Standardized Matrix Formulation

- The  $\tilde{X}^T \tilde{X}$  and  $\tilde{X}^T \tilde{Y}$  are correlation matrices

$$\tilde{X}^T \tilde{X} = \begin{bmatrix} 1 & r_{12} & \cdots & r_{1,p-1} \\ r_{12} & 1 & \cdots & r_{2,p-1} \\ \vdots & \vdots & \ddots & \vdots \\ r_{1n} & r_{2n} & \cdots & 1 \end{bmatrix} \quad \tilde{X}^T \tilde{Y} = \begin{bmatrix} r_{1y} \\ \vdots \\ r_{ny} \end{bmatrix}$$

- where  $r_{ij}$  is the correlation between the i<sup>th</sup> and j<sup>th</sup> regressors,  $r_{iy}$  is the correlation between the i<sup>th</sup> regressor and y

© Chris Mack, 2016 Data to Decisions 6

THE UNIVERSITY OF TEXAS  
AT AUSTIN

WHAT STARTS HERE CHANGES THE WORLD

### Consider a Two Regressor Model

$$\tilde{y} = \tilde{\beta}_1 \tilde{x}_1 + \tilde{\beta}_2 \tilde{x}_2 + \varepsilon \quad (\tilde{X}^T \tilde{X})^{-1} = \begin{bmatrix} \frac{1}{(1-r_{12}^2)} & \frac{-r_{12}}{(1-r_{12}^2)} \\ \frac{-r_{12}}{(1-r_{12}^2)} & \frac{1}{(1-r_{12}^2)} \end{bmatrix}$$

$$R^2 = \frac{r_{1y}^2 + r_{2y}^2 - 2r_{12}r_{1y}r_{2y}}{(1-r_{12}^2)} = \tilde{\beta}_1 r_{1y} + \tilde{\beta}_2 r_{2y} \quad \text{cov}(\tilde{\beta}_1, \tilde{\beta}_2) = \frac{-r_{12}S_\varepsilon^2}{(1-r_{12}^2)}$$

$$\tilde{\beta}_1 = \frac{r_{1y} - r_{12}r_{2y}}{(1-r_{12}^2)} \quad \tilde{\beta}_2 = \frac{r_{2y} - r_{12}r_{1y}}{(1-r_{12}^2)} \quad SE(\tilde{\beta}_1) = SE(\tilde{\beta}_2) = \frac{S_\varepsilon}{\sqrt{1-r_{12}^2}}$$

$$\tilde{h}_{ii} = h_{ii} = \frac{1}{n} + \frac{1}{(1-r_{12}^2)}(\tilde{x}_{1i}^2 + \tilde{x}_{2i}^2 - 2r_{12}\tilde{x}_{1i}\tilde{x}_{2i})$$

© Chris Mack, 2016 Data to Decisions 7

THE UNIVERSITY OF TEXAS  
AT AUSTIN

WHAT STARTS HERE CHANGES THE WORLD

### Effects of Standardization

- Standardizing the variables can sometimes be useful
  - Reduces correlations between  $\tilde{x}$  and  $\tilde{x}^2$ , and between interactions and primary terms
  - Less numerical instabilities for matrix inversion and OLS solution
- Standardization will not help with multicollinearity
  - We will use standardization later, with principle component analysis (PCA)

© Chris Mack, 2016 Data to Decisions 8

THE UNIVERSITY OF TEXAS  
AT AUSTIN

WHAT STARTS HERE CHANGES THE WORLD

### Lecture 48: What have we learned?

- What is variable standardization, and why is it used?
- In general, will standardization help with problems of multicollinearity?

© Chris Mack, 2016 Data to Decisions 9